

# Linear-Complexity Exponentially-Consistent Tests for Universal Outlying Sequence Detection

Yuheng Bu    Shaofeng Zou    Venugopal V. Veeravalli  
 University of Illinois at Urbana-Champaign  
 Email: bu3@illinois.edu, szou3@illinois.edu, vvv@illinois.edu

February 8, 2017

## Abstract

We study a universal outlying sequence detection problem, in which there are  $M$  sequences of samples out of which a small subset of outliers need to be detected. A sequence is considered as an outlier if the observations therein are generated by a distribution different from those generating the observations in the majority of the sequences. In the universal setting, the goal is to identify all the outliers without any knowledge about the underlying generating distributions. In prior work, this problem was studied as a universal hypothesis testing problem, and a generalized likelihood (GL) test was constructed and its asymptotic performance characterized. In this paper, we propose a different class of tests for this problem based on distribution clustering. Such tests are shown to be exponentially consistent and their time complexity is linear in the total number of sequences, in contrast with the GL test, which has time complexity that is exponential in the number of outliers. Furthermore, our tests based on clustering are applicable to more general scenarios. For example, when both the typical and outlier distributions form clusters, the clustering based test is exponentially consistent, but the GL test is not even applicable.

## 1 Introduction

In this paper, we study a universal outlying sequence detection problem, where there are  $M$  sequences of samples out of which a small subset are outliers and need to be detected. Each sequence consists of independent and identically distributed (i.i.d.) discrete observations. It is assumed that the observations in the majority of the sequences are distributed according to typical distributions. A sequence is considered as an outlier if its distribution is different from the typical distributions. We are interested in the universal setting of the problem, where nothing is known about the outlier and typical distributions, except that the outlier distributions are different from the typical distributions. The goal is to design a test, which does not depend on the typical and outlier distributions, to best discern all the outliers.

Outlying sequence detection finds possible applications in many domains (e.g., [1–3]). The problem of outlying sequence detection was studied as a universal outlier hypothesis testing problem for discrete observations in [4] and continuous observations in [5]. In [4], the exponential consistency of the generalized likelihood (GL) test under various universal settings was established. When there are a known number of identically distributed outliers, the GL test was shown to be asymptotically optimal as  $M$  goes to infinity.

However, the high time complexity, which is exponential in the number of outliers  $T$ , is a major drawback of the GL test when  $M$  and  $T$  are large. In this paper, we propose tests based on distribution clustering [6] for various scenarios. Such tests are shown to be exponentially consistent, with time complexity that is linear in  $M$  and independent of  $T$ .

The intuition for the clustering based test is that the typical distributions are usually closer to each other than the outlier distributions. Therefore, the typical distributions (and also possibly the outlier distributions)

will naturally form a cluster. This implies that the outlying sequence detection problem is equivalent to clustering the distributions using KL divergence as the distance metric (see also [7]).

We note that our problem is different from the classical distribution clustering problem [6, 8–10]. In the distribution clustering problem, the goal is to find the cluster structure with the lowest cost (sum of distance functions of each point in the cluster to the center). However, for our problem, we are given a sequence of samples from each distribution rather than the actual underlying distribution itself. Therefore, we are interested in the statistical performance, i.e., the error probability, of our test. Previous studies on approximation algorithms for distribution clustering [8–10] only show that the cost corresponding to the cluster structure returned by the approximation is within a  $\log K$  factor of the minimal cost, where  $K$  is the number of clusters. And there are no results showing that the approximation algorithms will converge to the minimal cost. Therefore, their results cannot be directly applied to our problem to provide a performance guarantee in a probabilistic sense.

Our contributions in this paper are as follows: (1) in all cases where the GL test is exponentially consistent, we construct tests based on clustering that are also exponentially consistent and have time complexity linear in  $M$ ; (2) we show that the tests based on clustering are applicable to more general scenarios. For example, when both the typical and outlier distributions form clusters, the clustering based test is exponentially consistent, but the GL test is not even applicable.

The rest of the paper is organized as follows. In Section 2, we describe the problem model. In Section 3, we introduce the GL test studied in [4]. In Section 4, we reformulate the outlying sequence detection problem as a distribution clustering problem. In Section 5, we propose linear complexity tests based on the K-means clustering algorithm. In Section 6, we provide some numerical results.

## 2 Problem Model

Throughout the paper, all distributions are defined on the finite set  $\mathcal{Y}$ , and  $\mathcal{P}(\mathcal{Y})$  denotes the set of all probability mass functions (pmfs) on  $\mathcal{Y}$ . All logarithms are with respect to the natural base.

We consider an outlying sequence detection problem, where there are in total  $M \geq 3$  data sequences denoted by  $Y^{(i)}$  for  $i = 1, \dots, M$ . Each data sequence  $Y^{(i)}$  consists of  $n$  i.i.d. samples  $Y_1^{(i)}, \dots, Y_n^{(i)}$ . The majority of the sequences are distributed according to typical distributions except for a subset  $S$  of outlying sequences, where  $S \subset \{1, \dots, M\}$  and  $1 \leq |S| = T < \frac{M}{2}$ . Each typical sequence  $j$  is distributed according to a typical distribution  $\pi_j \in \mathcal{P}(\mathcal{Y})$ ,  $j \notin S$ . Each outlying sequence  $i$  is distributed according to an outlier distribution  $\mu_i \in \mathcal{P}(\mathcal{Y})$ ,  $i \in S$ . Nothing is known about  $\mu_i$  and  $\pi_j$  except that  $\mu_i \neq \pi_j$ ,  $\forall i \in S, j \notin S$ ,  $S \subset \{1, \dots, M\}$ , and all of them have full support over a finite alphabet  $\mathcal{Y}$ .

We first study the case that all typical distributions are identical, i.e.,  $\pi_j = \pi$ ,  $\forall j \notin S$ .

We then study the case where both the outlier distributions  $\{\mu_i\}_{i \in S}$  and the typical distributions  $\{\pi_j\}_{j \notin S}$  are distinct. Moreover, the typical distributions and the outlier distributions form two clusters. More specifically,

$$\begin{aligned} \max_{i,j \in S} D(\mu_i \| \mu_j) &< \min_{i \in S, j \notin S} \{D(\mu_i \| \pi_j), D(\pi_j \| \mu_i)\}, \\ \max_{i,j \notin S} D(\pi_i \| \pi_j) &< \min_{i \in S, j \notin S} \{D(\mu_i \| \pi_j), D(\pi_j \| \mu_i)\}. \end{aligned} \quad (1)$$

This condition means that the divergence within the same cluster is less than the divergence between different clusters.

We use the notation  $\mathbf{y}^{(i)} = (y_1^{(i)}, \dots, y_n^{(i)})$ , where  $y_k^{(i)} \in \mathcal{Y}$  denotes the  $k$ -th observation of the  $i$ -th sequence. Let  $S$  be the set comprising all possible outlier subsets. For the hypothesis corresponding to an outlier subset

$S \in \mathcal{S}$ , the joint distribution of all the observations is given by

$$\begin{aligned} p_S(y^{Mn}) &= L_S(y^{Mn}, \{\mu_i\}_{i \in S}, \{\pi_j\}_{j \notin S}) \\ &= \prod_{k=1}^n \left\{ \prod_{i \in S} \mu_i(y_k^{(i)}) \prod_{j \notin S} \pi_j(y_k^{(j)}) \right\}, \end{aligned} \quad (2)$$

where  $L_S(y^{Mn}, \{\mu_i\}_{i \in S}, \{\pi_j\}_{j \notin S})$  denotes the likelihood, which is a function of the observations  $y^{Mn}$ ,  $\{\mu_i\}_{i=1}^M$  and  $\{\pi_j\}_{j=1}^M$ .

Our goal is to build distribution-free tests to detect the outlying sequences. The test can be captured by a universal rule  $\delta : \mathcal{Y}^{Mn} \rightarrow \mathcal{S}$ , which must not depend on  $\{\mu_i\}_{i=1}^M$  and  $\{\pi_j\}_{j=1}^M$ .

The performance of a universal test is gauged by the maximal probability of error, which is defined as

$$e(\delta, \{\mu_i\}_{i=1}^M, \{\pi_j\}_{j=1}^M) \triangleq \max_{S \in \mathcal{S}} \sum_{y^{Mn} : \delta(y^{Mn}) \neq S} p_S(y^{Mn}),$$

and the corresponding error exponent is defined as

$$\alpha(\delta, \{\mu_i\}_{i=1}^M, \{\pi_j\}_{j=1}^M) \triangleq \lim_{n \rightarrow \infty} -\frac{1}{n} \log e(\delta, \{\mu_i\}_{i=1}^M, \{\pi_j\}_{j=1}^M).$$

A universal test  $\delta$  is termed *universally exponentially consistent* if

$$\alpha(\delta, \{\mu_i\}_{i=1}^M, \{\pi_j\}_{j=1}^M) > 0,$$

for  $\mu_i \neq \pi_j$ ,  $i, j = 1, \dots, M$ .

### 3 Generalized Likelihood Test

In this section, we consider the case where the typical distributions are identical, i.e.,  $\pi_j = \pi$ ,  $\forall j \notin S$ . Let  $\gamma_i$  denote the empirical distribution of  $\mathbf{y}^{(i)}$ , and is defined as  $\gamma_i(y) \triangleq \frac{1}{m} |\{k = 1, \dots, m : y_k = y\}|$ , for each  $y \in \mathcal{Y}$ .

In the universal setting with  $\pi$  and  $\{\mu_i\}_{i \in S}$  unknown, conditioned on the set of outliers being  $S \in \mathcal{S}$ , we compute the generalized likelihood of  $y^{Mn}$  by replacing  $\pi$  and  $\{\mu_i\}_{i \in S}$  in (2) with their maximum likelihood estimates (MLEs)  $\{\hat{\mu}_i\}_{i \in S}$ , and  $\hat{\pi}_S$ , as

$$\hat{p}_S^{\text{univ}} = \hat{L}_S(y^{Mn}, \{\hat{\mu}_i\}_{i \in S}, \hat{\pi}_S). \quad (3)$$

The GL test then selects the hypothesis under which the GL is maximized (ties are broken arbitrarily), i.e.,

$$\delta_{\text{GL}}(y^{Mn}) = \arg \max_{S \in \mathcal{S}} \hat{p}_S^{\text{univ}}. \quad (4)$$

#### 3.1 Known Number of Outliers

We first consider the case in which the number of outliers, denoted by  $T \geq 1$ , is known at the outset, i.e.,  $\mathcal{S} = \{S : S \subset \{1, \dots, M\}, |S| = T\}$ .

We compute the generalized likelihood of  $y^{Mn}$  by replacing the  $\mu_i$ ,  $i \in S$  and  $\pi$  in (2) with their MLEs:  $\hat{\mu}_i \triangleq \gamma_i$ , and  $\hat{\pi}_S \triangleq \frac{\sum_{j \notin S} \gamma_j}{M - T}$ . Then the GL test in (4) is equivalent to

$$\delta_{\text{GL}}(y^{Mn}) = \arg \min_{S \subset \mathcal{S}} \sum_{j \notin S} D\left(\gamma_j \parallel \frac{\sum_{j \notin S} \gamma_j}{M - T}\right), \quad (5)$$

where  $D(p\|q)$  denotes the KL divergence between two distributions  $p, q \in \mathcal{P}(\mathcal{Y})$ , defined as

$$D(p\|q) \triangleq \sum_{y \in \mathcal{Y}} p(y) \log \left( \frac{p(y)}{q(y)} \right).$$

**Proposition 1.** [4, Theorem 9, 10] *When the number of outliers is known, the GL test in (5) is universally exponentially consistent. As  $M \rightarrow \infty$ , the achievable error exponent converges as*

$$\lim_{M \rightarrow \infty} \alpha \left( \delta_{\text{GL}}, \{\mu_i\}_{i=1}^M, \pi \right) = \lim_{M \rightarrow \infty} \min_{i=1, \dots, M} 2B(\mu_i, \pi),$$

where  $B(p, q)$  denotes the Bhattacharyya distance between two distributions  $p, q \in \mathcal{P}(\mathcal{Y})$ , defined as

$$B(p, q) \triangleq -\log \left( \sum_{y \in \mathcal{Y}} p(y)^{1/2} q(y)^{1/2} \right).$$

When all outlier sequences are identically distributed, i.e.,  $\mu_i = \mu \neq \pi$ ,  $i = 1, \dots, M$ , the achievable error exponent of the GL test in (5) converges to the optimal one achievable when both  $\mu$  and  $\pi$  are known.

Note that the number of hypotheses in the test (5) is  $\binom{M}{T}$ . Thus, exhaustive search over all possible hypotheses has time complexity that is polynomial in  $M$  and exponential in  $T$ .

### 3.2 Unknown Number of Identical Outliers

In this subsection, we consider the case where the number of outliers is unknown, i.e.,  $\mathcal{S} = \{S : S \subset \{1, \dots, M\}, 1 \leq |S| < M/2\}$ . Moreover, we assume that the outliers are identically distributed.

By replacing the  $\mu_i$ ,  $i \in S$ , and  $\pi$  in (2) with their MLEs  $\hat{\mu}_S = \hat{\mu}_i \triangleq \frac{\sum_{i \in S} \gamma_i}{|S|}$ , and  $\hat{\pi}_S \triangleq \frac{\sum_{j \notin S} \gamma_j}{M - |S|}$ , the GL test in (4) is equivalent to

$$\delta_{\text{GL}}(y^{Mn}) = \arg \min_{S \subset \mathcal{S}} \sum_{j \notin S} D \left( \gamma_j \parallel \frac{\sum_{j \notin S} \gamma_j}{M - |S|} \right) + \sum_{i \in S} D \left( \gamma_i \parallel \frac{\sum_{i \in S} \gamma_i}{|S|} \right). \quad (6)$$

**Proposition 2.** [4, Theorem 11] *When the number of the outliers is unknown, and  $1 \leq |S| < \frac{M}{2}$ , the GL test in (6) is universally exponentially consistent.*

Note that the number of hypotheses in the GL test (6) is  $\sum_{i=1}^{\lfloor M/2 \rfloor} \binom{M}{i}$ , which is exponential in  $M$ . The complexity of the test in (6) is even larger than that of the test in (5).

We also note that when the number of the outliers is unknown and outliers can be distinctly distributed, there cannot exist a universally exponentially consistent test [4, Theorem 12].

## 4 Problem Reformulation as Distribution Clustering

In order to construct low time complexity algorithms without sacrificing much in performance, we reformulate the outlying sequence detection problem as a distribution clustering problem. Although the distribution clustering problem is known to be NP-hard [8], there exists many approximation approaches, e.g. K-means algorithm [11], with time complexity that is linear in  $M$  and linear in the number of clusters.

The typical distributions are closer to each other than to any of the outlier distributions, and the same also holds for the outlier distributions when the outliers form a cluster. The outliers can be identified by clustering the empirical distributions into two clusters, where the cluster with more members contains all typical sequences, and the other cluster contains outliers.

More specifically, if we define the following cost function for distribution clustering

$$TC = \sum_{k=1}^K \sum_{i \in C^k} D(\gamma_i \| c^k), \quad (7)$$

where  $K$  is the number of clusters,  $c = \{c^1, \dots, c^K\}$  are the clustering centers, and disjoint partitions  $C = \{C^1, \dots, C^K\}$  denote the cluster assignment. As shown in [6, Proposition 1], for a given cluster assignment  $\{C^k\}_{k=1}^K$ , the total cost is minimized when  $c^k = \frac{\sum_{i \in C^k} p^i}{|C^k|}$ , which is the average of the distributions within the  $k$ -th cluster.

Then the GL test in (6) can be interpreted as a distribution clustering algorithm with  $K = 2$ . The first term in (6) corresponds to the cost in the typical cluster and the second term is the cost within the outlier cluster. The GL test decides on the cluster assignment that minimizes the cost among all possible cluster assignments.

For the case where the typical distributions are identically distributed, but outliers are not, it suffices to only cluster the empirical distributions of all typical sequences as shown in the GL test (5).

Thus, both the GL test in (5) and (6) are equivalent to empirical distribution clustering on the probability simplex using KL divergence as the distance metric.

While the distribution clustering problem itself is known to be NP-hard [8], there are many existing approximation algorithms with low complexity [11]. Here, we introduce the K-means distribution clustering algorithm in [6].

---

**Algorithm 1** K-means distribution clustering algorithm

---

**Input:**  $M$  distributions  $p^1, \dots, p^M$  defined on  $\mathcal{Y}$ , number of clusters  $K$ .

**Output:** partition set  $\{C^k\}_{k=1}^K$ .

**Initialization:**  $\{c_k\}_{k=1}^K$  (Specified in Algorithm 2 and 3)

**Method:**

**while** not convergence **do**

    {Assignment Step}

    Set  $C^k \leftarrow \emptyset$ ,  $1 \leq k \leq K$

**for**  $i = 1$  to  $M$  **do**

$C^k \leftarrow C^k \cup \{p^i\}$

        where  $k = \arg \min_k D(p^i \| c^k)$

**end for**

    {Re-estimation Step}

**for**  $k = 1$  to  $K$  **do**

$c^k \leftarrow \frac{\sum_{i \in C^k} p^i}{|C^k|}$

**end for**

**end while**

**Return**  $\{C^k\}_{k=1}^K$

---

For Algorithm 1, [6] only shows that the cost function in (7) is monotonically decreasing and it terminates in a finite number of steps at a partition that is locally optimal, i.e., the total loss cannot be decreased by either (a) the assignment step or by (b) changing the means of any existing clusters.

## 5 Clustering Based Tests

In this section, we propose linear complexity tests based on the K-means clustering algorithm. We show in all cases where the GL test is exponentially consistent, the clustering based tests using KL divergence as the distance metric are also exponentially consistent, while only taking linear time in  $M$ . For the case that the

typical and outlier distributions form two clusters, we show that the clustering based test is exponentially consistent, but the GL test is not even applicable.

## 5.1 Known Number of Outliers

We first consider the case where the number of outliers  $T$  is known and the typical distributions are identical. Algorithm 1 cannot be directly applied here, because the outlier distributions may not form a cluster and Algorithm 1 does not employ the knowledge of  $T$ . Motivated by the test in (5), we design Algorithm 2. The novelty of this algorithm lies in the construction of the first clustering center for the typical distribution and an iterative approach based on K-means to update it.

---

### Algorithm 2 Clustering with known number of outliers

---

**Input:**  $\gamma_1, \dots, \gamma_M$ , number of the outliers  $T$ .  
**Output:** A set of outliers  $S$ .  
**Initialization:**  
Choose one distribution  $\gamma^{(0)}$  from  $\gamma_1, \dots, \gamma_M$  arbitrarily  
**for**  $i = 1$  to  $M$  **do**  
    Compute  $D(\gamma_i \| \gamma^{(0)})$   
**end for**  
 $\hat{\pi} \leftarrow \gamma^*$   
where  $D(\gamma^* \| \gamma^{(0)})$  is the  $\lceil \frac{M}{2} \rceil$ -th element among all  $D(\gamma_i \| \gamma^{(0)})$   
**Method:**  
**while** not convergence **do**  
    {Assignment Step}  
Set  $S \leftarrow S^*$ ,  
where  $S^* = \arg \max_{S' \in \mathcal{S}, |S'|=T} \sum_{i \in S'} D(\gamma_i \| \hat{\pi})$   
    {Re-estimation Step}  
 $\hat{\pi} \leftarrow \frac{\sum_{j \notin S} \gamma_j}{M-T}$   
**end while**  
**Return**  $S$

---

Using the initialization in Algorithm 2,  $\gamma^*$  is generated from  $\pi$  with high probability. The intuition behind this is that: if  $\gamma^{(0)}$  is generated from typical distribution  $\pi$ , then only  $|S| < \frac{M}{2}$  empirical distributions which are generated from  $\mu_i$  are far from  $\gamma^{(0)}$ ; if  $\gamma^{(0)}$  is generated from some  $\mu_i$ , then there are at least  $M - |S| > \frac{M}{2}$  of  $D(\gamma_i \| \gamma^{(0)})$  concentrating at  $D(\pi \| \mu_i)$ . Thus the  $\lceil \frac{M}{2} \rceil$ -th element among  $D(\gamma_i \| \gamma^{(0)})$  is close to  $D(\pi \| \mu_i)$ , and  $\gamma^*$  is generated from  $\pi$  with high probability.

Let  $\delta_{c2}$  denote the test described in Algorithm 2, and  $\delta_{c2}^{(\ell)}$  denote the test that runs  $\ell$  iterations in Algorithm 2, where  $\ell = 1, 2, \dots$  is the number of iterations.

In the following theorem, we show that the test  $\delta_{c2}^{(1)}$  with only one iteration step is universally exponentially consistent.

**Theorem 1.** *For each  $M \geq 3$ , when the number of outliers  $T$  is known, the test  $\delta_{c2}^{(1)}$ , which runs one K-means iteration in Algorithm 2 is universally exponentially consistent. The achievable error exponent of  $\delta_{c2}^{(1)}$  can be upper bounded by*

$$\alpha\left(\delta_{c2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) < \lim_{M \rightarrow \infty} \min_{i=1, \dots, M} 2B(\mu_i, \pi). \quad (8)$$

Furthermore, the time complexity of the test  $\delta_{c2}^{(1)}$  is  $O(M)$ .

*Proof sketch.* Errors made by  $\delta_{c2}^{(1)}$  in the initialization step can be decomposed into two scenarios. If  $\gamma^{(0)}$  is generated from typical distribution  $\pi$ , an error occurs when  $\hat{\pi}$  is actually generated from an outlier

distribution. The probability of this event can be upper bounded by the probability of the following event  $E_1 = \{D(\gamma_i \parallel \gamma_{j_1}) < D(\gamma_{j_2} \parallel \gamma_{j_1}), \exists i \in S, \exists j_1, j_2 \notin S\}$ . If  $\gamma^{(0)}$  is generated from an outlier distribution, the error probability can be upper bounded by the probability of the following event  $E_2 = \{D(\gamma_{j_1} \parallel \gamma_{i_1}) < D(\gamma_{i_2} \parallel \gamma_{i_1}) < D(\gamma_{j_2} \parallel \gamma_{i_1}), \exists i_1, i_2 \in S, \exists j_1, j_2 \notin S\}$ . By Sanov's theorem, we can prove that the probability of  $E_1$  and  $E_2$  decay exponentially fast.

The error probability in the assignment step can be upper bounded by the probability of the same event  $E_1$ , which decays exponentially fast by Sanov's theorem.

Moreover, the assignment step in Algorithm 2 can be solved in linear time  $O(M)$  [12], which is independent of  $T$ .

The details of the proof can be found in Appendix B.  $\square$

Comparison of Proposition 1 and Theorem 1 shows that  $\delta_{c_2}^{(1)}$  has a smaller error exponent than that of the GL test in (5) as  $M \rightarrow \infty$ , but has a linear time complexity.

In the following theorem, we further show that the performance of Algorithm 2 will improve with more iterations.

**Theorem 2.** *For each  $M \geq 3$ , when the number of outliers  $T$  is known, the test  $\delta_{c_2}^{(\ell)}$  is universally exponentially consistent. As  $M \rightarrow \infty$ , the achievable error exponent of  $\delta_{c_2}^{(\ell)}$  in Algorithm 2 can be lower bounded by*

$$\lim_{M \rightarrow \infty} \alpha\left(\delta_{c_2}^{(\ell)}, \{\mu_i\}_{i=1}^M, \pi\right) \geq \alpha\left(\delta_{c_2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) \quad (9)$$

Furthermore, the time complexity of the test  $\delta_{c_2}^{(\ell)}$  is  $O(M\ell)$ .

*Proof sketch.* It is shown in Theorem 1 and Proposition 1, both the test  $\delta_{c_2}^{(1)}$  and the GL test  $\delta_{\text{GL}}$  are exponentially consistent, and the error exponent of  $\delta_{\text{GL}}$  is larger than that of  $\delta_{c_2}^{(1)}$ , when  $M \rightarrow \infty$ . Since  $\mathbb{P}(\delta_{c_2}^{(1)} \neq S) \leq \mathbb{P}(\delta_{c_2}^{(1)} \neq S \text{ or } \delta_{\text{GL}} \neq S) \leq \mathbb{P}(\delta_{c_2}^{(1)} \neq S) + \mathbb{P}(\delta_{\text{GL}} \neq S)$ , we can conclude that  $\mathbb{P}(\delta_{c_2}^{(1)} \neq S \text{ or } \delta_{\text{GL}} \neq S)$  also decays exponentially fast with the same error exponent of  $\delta_{c_2}^{(1)}$ , which means that  $\delta_{\text{GL}}$  and  $\delta_{c_2}^{(1)}$  both output the same correct  $S$  with high probability. Given that  $\delta_{c_2}^{(1)}$  and  $\delta_{\text{GL}}$  have same outcomes, i.e., the initialization achieves the global optimum of the cost function, then running  $\ell$  steps will not change the output of  $\delta_{c_2}^{(1)}$ . Thus,  $\delta_{c_2}^{(\ell)}$  at least achieves the error exponent of  $\delta_{c_2}^{(1)}$ . The details of the proof can be found in Appendix C.  $\square$

## 5.2 Unknown Number of Identical Outliers

In this section, we consider the case where the number of outliers is unknown and both typical and outlier distributions are identical. Motivated by the test in (6), we design the following initialization algorithm to set the clustering centers in Algorithm 1.

---

### Algorithm 3 Clustering with unknown number of outliers

---

**Input:**  $M$  empirical distributions  $\gamma_1, \dots, \gamma_M$  defined on finite alphabet  $\mathcal{Y}$ .

**Output:** A set of outliers  $S$ .

**Initialization:**

Choose one distribution  $\gamma^{(0)}$  arbitrarily,

$c^1 \leftarrow \arg \max_{\gamma_i} D(\gamma_i \parallel \gamma^{(0)})$

$c^2 \leftarrow \gamma^{(0)}$

**Method:** Same as in Algorithm 1 with  $K = 2$

**Return**  $C^1$  and  $C^2$

---

It can be seen that,  $c_1$  and  $c_2$  chosen by the initialization step in Algorithm 3 are generated by  $\pi$  and  $\mu$ , with high probability.

Let  $\delta_{c3}$  denote the test described in Algorithm 3, and  $\delta_{c3}^{(\ell)}$  denote the test that runs  $\ell$  iterations in Algorithm 2, where  $\ell = 1, 2, \dots$  is the number of iterations. The following theorem shows that the clustering based algorithm  $\delta_{c3}^{(\ell)}$ , is universally exponentially consistent, and has time complexity linear in  $M$ .

**Theorem 3.** *For each  $M \geq 3$ , when the number of outliers is unknown, the test  $\delta_{c3}^{(\ell)}$ , which runs  $\ell$  steps of Algorithm 3, is exponentially consistent, and has time complexity  $O(M\ell)$ .*

*Proof sketch.* The exponential consistency of  $\delta_{c3}^{(\ell)}$  can be established using similar techniques to those in Theorem 1 and Theorem 2. The major difference between the proof of Theorem 1 and Theorem 3 is that there are two clustering centers in the initialization step and assignment step in Algorithm 3. The details of the proof can be found in Appendix D.  $\square$

### 5.3 Typical and Outlier Distributions Forming Clusters

In this subsection, we consider the case that both the outlier distributions  $\{\mu_i\}_{i \in S}$  and the typical distributions  $\{\pi_j\}_{j \notin S}$  are distinct. Moreover, the typical distributions and the outlier distributions form clusters as defined in (1), which means that the divergence within the same cluster is always less than the divergence between different clusters.

The following theorem shows that under the condition (1), the one step test  $\delta_{c3}^{(1)}$  proposed in Algorithm 3 is universally exponentially consistent, and has time complexity linear in  $M$ .

**Theorem 4.** *For each  $M \geq 3$ , when both the outlier distributions  $\{\mu_i\}_{i \in S}$  and the typical distributions  $\{\pi_j\}_{j \notin S}$  form clusters, i.e. condition (1) holds, the test  $\delta_{c3}^{(1)}$ , which runs one step of Algorithm 3, is universally exponentially consistent, and has time complexity  $O(M)$ .*

*Proof sketch.* The exponential consistency of  $\delta_{c3}^{(1)}$  can be established using similar techniques as shown in Theorem 3. The details of the proof can be found in Appendix E.  $\square$

The GL approach of replacing the true distribution in (2) by their MLEs leads to identical likelihood estimates for each hypothesis. Thus, the GL approach is not applicable here. One could apply the test in (6) to this problem, but it can be shown (see Appendix F) that the test in (6) is not universally exponentially consistent, even if condition (1) holds.

## 6 Numerical Results

In this section, we compare the performance of the GL test  $\delta_{GL}$ , the clustering based tests  $\delta_{c2}$ ,  $\delta_{c3}$  and the one step tests  $\delta_{c2}^{(1)}$ ,  $\delta_{c3}^{(1)}$ . We consider the case with identical typical distribution. We set  $\pi$  to be the uniform distribution with alphabet size 10, and generate outlier distributions randomly.

We first simulate the case with distinct outliers where  $T$  is known. We choose  $M = 20$  sequences and  $T = 3$  outliers. In Fig. 1, we plot  $\log P_e$  as a function of  $n$  for  $\delta_{GL}$ ,  $\delta_{c2}$  and  $\delta_{c2}^{(1)}$ . As we can see from Fig. 1,  $\delta_{c2}^{(1)}$  and  $\delta_{c2}$  are both exponentially consistent, and  $\delta_{c2}$  outperforms  $\delta_{c2}^{(1)}$  as shown in Theorem 2. Comparison between  $\delta_{c2}$  and  $\delta_{GL}$  shows that without sacrificing much in performance,  $\delta_{c2}$  is about 50 times faster than  $\delta_{GL}$ .

We then simulate the case with unknown number of identical outliers. We set  $M = 100$  sequences and  $T = 10$  outliers. Fig. 2 shows that  $\delta_{c3}$  outperforms  $\delta_{c3}^{(1)}$ . We note that running the clustering based tests for 5000 times takes 5 minutes on a 3.6 GHz i7 CPU. However, the GL test is not feasible here, since the number of hypotheses needs to search over is exponential in  $M$ .



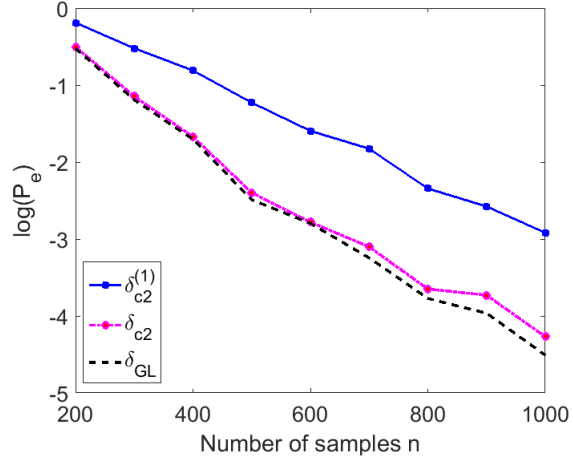


Figure 1: Comparison of test  $\delta_{c2}^{(1)}$ ,  $\delta_{c2}$ ,  $\delta_{GL}$ , known number of distinct outliers

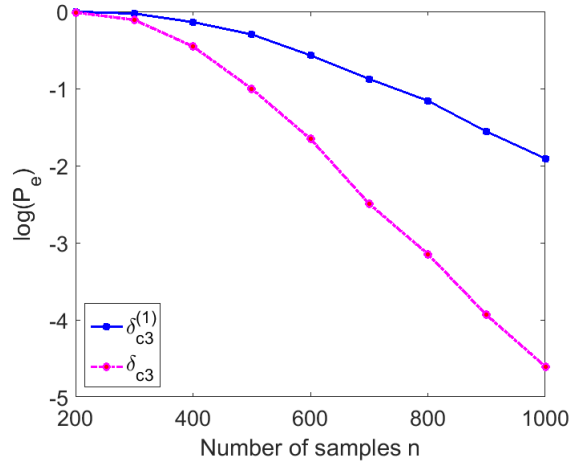


Figure 2: Comparison of test  $\delta_{c3}$ ,  $\delta_{c3}^{(1)}$ , unknown number of identical outliers

# Appendix

## A Useful Lemmas

Our proofs rely on the following lemmas provided in [13].

**Lemma 1.** [13, Lemma 1] Let  $\mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(J)}$  be mutually independent random vectors with each  $\mathbf{Y}^{(j)}$ ,  $j = 1, \dots, J$ , being  $n$  i.i.d. repetitions of a random variable distributed according to  $p_j \in \mathcal{P}(\mathcal{Y})$ . Let  $A_n$  be the set of all  $J$  tuples  $(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(J)}) \in \mathcal{Y}^{Jn}$  whose empirical distributions  $(\gamma_1, \dots, \gamma_J)$  lie in a closed set  $E \in \mathcal{P}(\mathcal{Y})^J$ . Then, it holds that

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P} \left\{ \left( \mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(J)} \right) \in A_n \right\} = \min_{(q_1, \dots, q_J) \in E} \sum_{j=1}^J D(q_j \| p_j). \quad (10)$$

**Lemma 2.** [13, Lemma 2] For any two pmfs  $p_1, p_2 \in \mathcal{P}(\mathcal{Y})$  with full supports, it holds that

$$2B(p_1, p_2) = \min_{q \in \mathcal{P}(\mathcal{Y})} \left( D(q \| p_1) + D(q \| p_2) \right). \quad (11)$$

In particular, the minimum on the right side is achieved by

$$q^* = \frac{p_1^{1/2}(y)p_1^{1/2}(y)}{\sum_{y \in \mathcal{Y}} p_1^{1/2}(y)p_1^{1/2}(y)}, \quad y \in \mathcal{Y}. \quad (12)$$

## B Proof of Theorem 1

Due to the structure of the test we know that errors occur at two different steps:

1. **Initialization Step:** The constructed clustering center for typical sequences  $\hat{\pi}$  is actually generated from the outlier distribution.
2. **Assignment Step:** Given correct clustering center  $\hat{\pi}$ , the empirical distribution of the outlying sequence is more close to the clustering center  $\hat{\pi}$ .

We first use the event  $E$  to denote that errors occur at the initialization step. Since we use the  $\gamma^*$  as the clustering center  $\hat{\pi}$ , where  $D(\gamma^* \| \gamma^{(0)})$  is the  $\lceil \frac{M}{2} \rceil$ -th element among all  $D(\gamma_i \| \gamma^{(0)})$ . It is difficult to write the explicit form of the event  $E$ . However, we can find upper bounds for the probability of  $E$  in the following two scenarios.

If  $\gamma^{(0)}$  is generated from typical distribution  $\pi$ , an error occurs when  $\hat{\pi}$  is actually generated from an outlier distribution, then  $D(\gamma_i \| \gamma^{(0)}) \leq D(\gamma_j \| \gamma^{(0)})$  must hold for some  $i \in S, j \notin S$ . Due to the arbitrariness of  $\gamma^{(0)}$ , the probability of this error event can be upper bounded by the probability of the following event

$$E_1 \triangleq \{D(\gamma_i \| \gamma_{j_1}) \leq D(\gamma_{j_2} \| \gamma_{j_1}), \exists i \in S, \exists j_1, j_2 \notin S\}. \quad (13)$$

If  $\gamma^{(0)}$  is generated from an outlier distribution, the error probability can be upper bounded by the probability of the following event

$$E_2 \triangleq \{D(\gamma_{j_1} \| \gamma_{i_1}) < D(\gamma_{i_2} \| \gamma_{i_1}) < D(\gamma_{j_2} \| \gamma_{i_1}), \exists i_1, i_2 \in S, \exists j_1, j_2 \notin S\}. \quad (14)$$

Thus,  $\mathbb{P}(E) \leq \mathbb{P}(E_1) + \mathbb{P}(E_2)$ .

We then use  $F$  to denote that errors occur at the assignment step, then

$$F \triangleq E^C \cap \{D(\gamma_i \|\hat{\pi}) \leq D(\gamma_j \|\hat{\pi}), \exists i \in S, \exists j \notin S\}. \quad (15)$$

Note that  $F \subset E_1$ , then the probability of error event  $F$  can be upper bounded by that of the event  $E_1$ .

The error probability of the test  $\delta_{c_2}^{(1)}$  can be bounded by

$$\mathbb{P}(F) \leq e\left(\delta_{c_2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) = \mathbb{P}(E \cup F) \leq \mathbb{P}(E) + \mathbb{P}(F). \quad (16)$$

The right hand side of (16) can be further bounded by

$$\begin{aligned} & \mathbb{P}(E) + \mathbb{P}(F) \\ & \leq \mathbb{P}(E_1) + \mathbb{P}(E_2) + \mathbb{P}(E_1) \\ & \leq (M - |S|)^2 |S|^2 \left( 2 \max_{i \in S} \mathbb{P}(D(\gamma_i \|\gamma_{j_1}) \leq D(\gamma_{j_2} \|\gamma_{j_1})) + \max_{i_1, i_2 \in S} \mathbb{P}(D(\gamma_{j_1} \|\gamma_{i_1}) < D(\gamma_{i_2} \|\gamma_{i_1}) < D(\gamma_{j_2} \|\gamma_{i_1})) \right), \end{aligned} \quad (17)$$

where  $j_1, j_2 \notin S$ .

As for the left hand side of (16), we have

$$\mathbb{P}(F) \geq \max_{i \in S} \mathbb{P}(\{D(\gamma_i \|\gamma_{j_1}) \leq D(\gamma_{j_2} \|\gamma_{j_1})\}), \quad (18)$$

where  $j_1, j_2 \notin S$ .

From Lemma 1, we know the exponent can be computed as

$$\begin{aligned} \lim_{n \rightarrow \infty} -\frac{1}{n} \log \max_{i \in S} \mathbb{P}(D(\gamma_i \|\gamma_{j_1}) \leq D(\gamma_{j_2} \|\gamma_{j_1})) &= \min_{q_1, q_2, q_3 \in C_1, i \in S} D(q_1 \|\mu_i) + D(q_2 \|\pi) + D(q_3 \|\pi) \triangleq \alpha_1 \\ C_1 &= \{(q_1, q_2, q_3) : D(q_1 \|\mu_i) \leq D(q_3 \|\mu_i)\}. \end{aligned} \quad (19)$$

$$\begin{aligned} & \lim_{n \rightarrow \infty} -\frac{1}{n} \log \max_{i_1, i_2 \in S} \mathbb{P}(D(\gamma_{j_1} \|\gamma_{i_1}) < D(\gamma_{i_2} \|\gamma_{i_1}) < D(\gamma_{j_2} \|\gamma_{i_1})) \\ &= \min_{q_1, q_2, q_3, q_4 \in C_2, i_1, i_2 \in S} D(q_1 \|\pi) + D(q_2 \|\pi) + D(q_3 \|\mu_{i_1}) + D(q_4 \|\mu_{i_2}) \triangleq \alpha_2 \\ & C_2 = \{(q_1, q_2, q_3, q_4) : D(q_1 \|\pi) < D(q_4 \|\pi) < D(q_2 \|\pi)\}. \end{aligned} \quad (20)$$

Due to the fact that the objective function of (19) can only be zero at a collection  $q_1 = \mu_i, q_2 = q_3 = \pi$ , which are not in the constraint sets. And the objective function of (20) can only achieve zero when  $q_1 = q_2 = \pi, q_3 = \mu_{i_1}, q_4 = \mu_{i_2}$ , which are not in the constraint sets, either. Thus, we can conclude that  $\alpha_1, \alpha_2 > 0$ . From the fact that  $\lim_{n \rightarrow \infty} \frac{\log M(M-|S|)}{n} = 0$ , and combining with (17) and (18), we get that

$$\alpha_1 \geq \alpha\left(\delta_{c_2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) = \lim_{n \rightarrow \infty} -\frac{1}{n} \log e\left(\delta_{c_2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) \geq \alpha_1 \wedge \alpha_2. \quad (21)$$

This result shows that the one step test  $\delta_{c_2}^{(1)}$  is universally exponentially consistent. We note that

$$\alpha_1 = \min_{q_1, q_2, q_3 \in C_1, i \in S} D(q_1 \|\mu_i) + D(q_2 \|\pi) + D(q_3 \|\pi), \quad (22)$$

where  $C_1 = \{D(q_1 \|\mu_i) \leq D(q_3 \|\mu_i)\}$ . If we add the constrain that  $q_1 = q_3$ , i.e.,  $C'_1 = \{D(q_1 \|\mu_i) \leq D(q_3 \|\mu_i), q_1 = q_3\}$ , then  $C'_1 \subset C_1$ , and  $D(q \|\mu_i) \leq D(q \|\mu_i)$  holds, thus

$$\begin{aligned} \alpha_1 &< \min_{q_1, q_2, q_3 \in C'_1, i \in S} D(q_1 \|\mu_i) + D(q_2 \|\pi) + D(q_3 \|\pi) \\ &= \min_{q, q_2 \in C'_1, i \in S} D(q \|\mu_i) + D(q_2 \|\pi) + D(q \|\pi) \\ &= \min_{q \in \mathcal{P}(\mathcal{Y}), i \in S} D(q \|\mu_i) + D(q \|\pi). \end{aligned} \quad (23)$$

Here, we can set  $q_2 = \pi$ , due to our relaxation on set  $C'_1$ . From Lemma 2, we know that the minimum is the Bhattacharyya distance between the distribution  $\mu_i$  and  $\pi$

$$\alpha_1 < \min_{q \in \mathcal{P}(\mathcal{Y}), i \in S} D(q \parallel \mu_i) + D(q \parallel \pi) = \min_{i \in S} B(\mu_i, \pi). \quad (24)$$

Thus, as  $M \rightarrow \infty$ , we have

$$\alpha\left(\delta_{c2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) \leq \alpha_1 < \lim_{M \rightarrow \infty} \min_{i \in S} B(\mu_i, \pi). \quad (25)$$

As for the time complexity, it is obvious that the initialization step in Algorithm 2 can be executed within  $O(M)$  time. And the assignment step in Algorithm 2, which finds the largest  $T$  elements from size  $M$  array, can be solved in linear time  $O(M)$  using the algorithm proposed in [12]. Thus the overall time complexity is linear in  $M$  and independent of  $T$ .

Comparison of Proposition 1 and Theorem 1 shows that  $\delta_{c2}^{(1)}$  has a smaller error exponent than that of the GL test in (5) as  $M \rightarrow \infty$ , but has a linear time complexity.

## C Proof of Theorem 2

From Theorem 1 and Proposition 1, we know that both the one-step test  $\delta_{c2}^{(1)}$  and the GL test  $\delta_{GL}$  are exponentially consistent. We denote the error exponent of one-step K-means test as  $\alpha_{c2}^{(1)}$  and the error exponent of GL test as  $\alpha_{GL}$ . Then we have

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P}(\delta_{c2}^{(1)}(y^{Mn}) \neq S) = \alpha_{c2}^{(1)}, \quad (26)$$

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P}(\delta_{GL}(y^{Mn}) \neq S) = \alpha_{GL}. \quad (27)$$

Let

$$A \triangleq \{y^{Mn} : \delta_{c2}^{(1)}(y^{Mn}) \neq S\}, \quad B \triangleq \{y^{Mn} : \delta_{GL}(y^{Mn}) \neq S\}. \quad (28)$$

Then the set  $\{y^{Mn} : \delta_{c2}^{(1)}(y^{Mn}) = \delta_{GL}(y^{Mn}) = S\}$ , represented as  $A^C \cap B^C$ , can be further lower bounded by

$$\begin{aligned} \mathbb{P}(A^C \cap B^C) &= \mathbb{P}(A^C) + \mathbb{P}(B^C) - \mathbb{P}(A^C \cup B^C) \\ &\geq 1 - c_1 e^{-n\alpha_{c2}^{(1)}} - c_2 e^{-n\alpha_{GL}}, \end{aligned} \quad (29)$$

where  $c_1$  and  $c_2$  are some constant independent of  $n$ . Thus,

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log(1 - \mathbb{P}(A^C \cap B^C)) = \alpha_{c2}^{(1)} \wedge \alpha_{GL}. \quad (30)$$

This shows that the error probability decays exponentially fast. The one-step test  $\delta_{c2}^{(1)}$  and the GL test  $\delta_{GL}$  output the same correct  $S$  with high probability. Note that under the condition that  $\delta_{c2}^{(1)}$  and  $\delta_{GL}$  give the same outcome, running more iterations in Algorithm 2 will not change the output, thus

$$e\left(\delta_{c2}^{(\ell)}, \{\mu_i\}_{i=1}^M, \pi\right) \leq 1 - \mathbb{P}(A^C \cap B^C). \quad (31)$$

The error exponent of running  $\ell$  iterations will be lower bounded by

$$\alpha\left(\delta_{c2}^{(\ell)}, \{\mu_i\}_{i=1}^M, \pi\right) = \lim_{n \rightarrow \infty} -\frac{1}{n} \log e\left(\delta_{c2}^{(\ell)}, \{\mu_i\}_{i=1}^M, \pi\right) \geq \lim_{n \rightarrow \infty} -\frac{1}{n} \log(1 - \mathbb{P}(A^C \cap B^C)) = \alpha_{c2}^{(1)} \wedge \alpha_{GL}. \quad (32)$$

As  $M \rightarrow \infty$ , we know that  $\alpha_{c2}^{(1)} < \alpha_{GL}$  from Theorem 1, then we can conclude that

$$\lim_{M \rightarrow \infty} \alpha\left(\delta_{c2}^{(\ell)}, \{\mu_i\}_{i=1}^M, \pi\right) \geq \lim_{M \rightarrow \infty} \alpha\left(\delta_{GL}, \{\mu_i\}_{i=1}^M, \pi\right) \wedge \alpha\left(\delta_{c2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right) = \alpha\left(\delta_{c2}^{(1)}, \{\mu_i\}_{i=1}^M, \pi\right). \quad (33)$$

As for the time complexity, since each iteration has the complexity  $O(M)$ ,  $\delta_{c2}^{(\ell)}$  which runs  $\ell$  steps of iteration has complexity  $O(M\ell)$ .

## D Proof of Theorem 3

The exponential consistency of  $\delta_{c3}^{(\ell)}$  can be established using similar techniques to those in Theorem 1 and Theorem 2. The major difference between the proof of Theorem 1 and Theorem 3 is that there are two clustering centers in the initialization step and assignment step in Algorithm 3.

We first establish the exponential consistency of the one-step test  $\delta_{c3}^{(1)}$ .

Due to the structure of the test we know that errors occur at two different steps:

1. **Initialization Step:** The constructed clustering center for typical sequences  $\hat{\pi}$  and outliers  $\hat{\mu}$  are actually coming from the same distribution.
2. **Assignment Step:** The empirical distribution of the outlying sequence is more close to the clustering center of the typical sequence  $\hat{\pi}$ , and vice versa.

Note that in Algorithm 3,  $\gamma^{(0)}$  is chosen arbitrarily, so  $\gamma^{(0)}$  can be generated from  $\pi$  or  $\mu$ . We use  $\hat{\pi}$  and  $\hat{\mu}$  to denote that  $\gamma^{(0)}$  is generated from  $\pi$  or  $\mu$ , respectively. We first use the event  $E$  to denote that errors occur at the initialization step, the error event  $E$  can be decomposed as two parts:

$$E \triangleq \left\{ \max_{j \notin S} D(\gamma_j \| \hat{\pi}) > \max_{i \in S} D(\gamma_i \| \hat{\pi}) \right\} \cup \left\{ \max_{i \in S} D(\gamma_i \| \hat{\mu}) > \max_{j \notin S} D(\gamma_j \| \hat{\mu}) \right\}. \quad (34)$$

Denote

$$A_i \triangleq \left\{ \max_{j \notin S} D(\gamma_j \| \hat{\pi}) > D(\gamma_i \| \hat{\pi}) \right\}, \quad \forall i \in S, \quad (35)$$

$$B_j \triangleq \left\{ \max_{i \in S} D(\gamma_i \| \hat{\mu}) > D(\gamma_j \| \hat{\mu}) \right\}, \quad \forall j \notin S. \quad (36)$$

Then,

$$E = \left( \bigcap_{i \in S} A_i \right) \cup \left( \bigcap_{j \notin S} B_j \right). \quad (37)$$

The error event at the assignment step can be written as

$$F \triangleq E^C \cap \left\{ \max_{j \notin S} D(\gamma_j \| \hat{\pi}) - D(\gamma_j \| \hat{\mu}) > 0 \right\} \cup \left\{ \max_{i \in S} D(\gamma_i \| \hat{\mu}) - D(\gamma_i \| \hat{\pi}) > 0 \right\} = F_1 \cup F_2, \quad (38)$$

where

$$F_1 \triangleq E^C \cap \left\{ \max_{j \notin S} D(\gamma_j \| \hat{\pi}) - D(\gamma_j \| \hat{\mu}) > 0 \right\}, \quad (39)$$

$$F_2 \triangleq E^C \cap \left\{ \max_{i \in S} D(\gamma_i \| \hat{\mu}) - D(\gamma_i \| \hat{\pi}) > 0 \right\}. \quad (40)$$

Thus, we can upper bound the error probability of the one-step test  $\delta_{c3}^{(1)}$  by

$$e\left(\delta_{c3}^{(1)}, \mu, \pi\right) = \mathbb{P}(E \cup F) \leq \mathbb{P}(E) + \mathbb{P}(F). \quad (41)$$

$$\begin{aligned} \mathbb{P}(E) &\leq \mathbb{P}\left(\bigcap_{i \in S} A_i\right) + \mathbb{P}\left(\bigcap_{j \notin S} B_j\right) \leq \mathbb{P}(A_i) + \mathbb{P}(B_j) \\ &\leq (M - |S|)\mathbb{P}(D(\gamma_j \| \hat{\pi}) > D(\gamma_i \| \hat{\pi})) + |S|\mathbb{P}(D(\gamma_i \| \hat{\mu}) > D(\gamma_j \| \hat{\mu})) \\ &\leq (M - |S|)^2\mathbb{P}(D(\gamma_{j_1} \| \gamma_{j_2}) > D(\gamma_i \| \gamma_{j_2})) + |S|^2\mathbb{P}(D(\gamma_{i_1} \| \gamma_{i_2}) > D(\gamma_j \| \gamma_{i_2})) \end{aligned} \quad (42)$$

for  $j, j_1, j_2 \notin S$  and  $i, i_1, i_2 \in S$ . From lemma 1, the exponent can be computed as

$$\begin{aligned} \lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P}(D(\gamma_{j_1} \parallel \gamma_{j_2}) > D(\gamma_i \parallel \gamma_{j_2})) &= \min_{q_1, q_2, q_3 \in C_3} D(q_1 \parallel \pi) + D(q_2 \parallel \pi) + D(q_3 \parallel \mu) \triangleq \alpha_3 \\ C_3 &= \{(q_1, q_2, q_3) : D(q_1 \parallel q_2) > D(q_3 \parallel q_2)\}, \end{aligned} \quad (43)$$

$$\begin{aligned} \lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P}(D(\gamma_{i_1} \parallel \gamma_{i_2}) > D(\gamma_j \parallel \gamma_{i_2})) &= \min_{q_1, q_2, q_3 \in C_4} D(q_1 \parallel \mu) + D(q_2 \parallel \mu) + D(q_3 \parallel \pi) \triangleq \alpha_4 \\ C_4 &= \{(q_1, q_2, q_3) : D(q_1 \parallel q_2) > D(q_3 \parallel q_2)\}. \end{aligned} \quad (44)$$

We can upper bound  $\mathbb{P}(F)$  by union bounds,

$$\begin{aligned} \mathbb{P}(F) &\leq \mathbb{P}(F_1) + \mathbb{P}(F_2) \\ &\leq \mathbb{P}\left(\bigcup_{j \notin S} \{D(\gamma_j \parallel \hat{\pi}) - D(\gamma_j \parallel \hat{\mu}) > 0\}\right) + \mathbb{P}\left(\bigcup_{i \in S} \{D(\gamma_i \parallel \hat{\mu}) - D(\gamma_i \parallel \hat{\pi}) > 0\}\right) \\ &\leq |S|(M - |S|)^2 \mathbb{P}(D(\gamma_{j_1} \parallel \gamma_{j_2}) - D(\gamma_{j_1} \parallel \gamma_i) > 0) + |S|^2(M - |S|) \mathbb{P}(D(\gamma_{i_1} \parallel \gamma_{i_2}) - D(\gamma_{i_1} \parallel \gamma_j) > 0), \end{aligned} \quad (45)$$

where  $j, j_1, j_2 \notin S$  and  $i, i_1, i_2 \in S$ . From lemma 1, the exponent can be computed as

$$\begin{aligned} \lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P}(D(\gamma_{j_1} \parallel \gamma_{j_2}) - D(\gamma_{j_1} \parallel \gamma_i) > 0) &= \min_{q_1, q_2, q_3 \in C_5} D(q_1 \parallel \pi) + D(q_2 \parallel \pi) + D(q_3 \parallel \mu) \triangleq \alpha_5 \\ C_5 &= \{(q_1, q_2, q_3) : D(q_1 \parallel q_2) > D(q_1 \parallel q_3)\}, \end{aligned} \quad (46)$$

$$\begin{aligned} \lim_{n \rightarrow \infty} -\frac{1}{n} \log \mathbb{P}(D(\gamma_{i_1} \parallel \gamma_{i_2}) - D(\gamma_{i_1} \parallel \gamma_j) > 0) &= \min_{q_1, q_2, q_3 \in C_6} D(q_1 \parallel \mu) + D(q_2 \parallel \mu) + D(q_3 \parallel \pi) \triangleq \alpha_6 \\ C_6 &= \{(q_1, q_2, q_3) : D(q_1 \parallel q_2) > D(q_1 \parallel q_3)\}. \end{aligned} \quad (47)$$

Due to the fact that the objective function of (43), (46) can only be zero at a collection  $q_1 = q_2 = \pi$ ,  $q_3 = \mu$ , which are not in the constraint sets. And the objective function of (44), (47) can only achieve zero when  $q_1 = q_2 = \mu$ ,  $q_3 = \pi$ , which are not in the constraint sets, either. Thus, we can conclude that  $\alpha_3, \alpha_4, \alpha_5, \alpha_6 > 0$ . From that  $\lim_{n \rightarrow \infty} \frac{\log M(M - |S|)}{n} = 0$ , we get that

$$\alpha\left(\delta_{c3}^{(1)}, \mu, \pi\right) = \lim_{n \rightarrow \infty} -\frac{1}{n} \log e\left(\delta_{c3}^{(1)}, \mu, \pi\right) \geq \min\{\alpha_3, \alpha_4, \alpha_5, \alpha_6\}. \quad (48)$$

From the above argument and Proposition 2, we know that both the one-step test  $\delta_{c3}^{(1)}$  and the GL test  $\delta_{\text{GL}}$  are exponentially consistent. Thus, based on the same technique used in the proof of Theorem 2, we can establish the exponential consistency of the test  $\delta_{c3}^{(\ell)}$  proposed in Algorithm 3.

As for the time complexity, since each iteration has the complexity  $O(M)$ ,  $\delta_{c3}^{(\ell)}$  which runs  $\ell$  steps of iteration has complexity  $O(M\ell)$ .

## E Proof of Theorem 4

The exponential consistency of  $\delta_{c3}^{(1)}$  for the case where typical and outlier distributions form clusters can be established using the same techniques as shown in Theorem 3. The major difference between the proof of Theorem 3 and Theorem 4 is that both the typical distributions and the outlier distributions are distinct.

Using the same events defined in Appendix D, the error probability of the one-step test  $\delta_{c3}^{(1)}$  can be upper bounded by

$$e\left(\delta_{c3}^{(1)}, \{\mu_i\}_{i=1}^M, \{\pi_i\}_{i=1}^M\right) = \mathbb{P}(E \cup F) \leq \mathbb{P}(E) + \mathbb{P}(F). \quad (49)$$

$$\begin{aligned}
\mathbb{P}(E) &\leq \mathbb{P}\left(\bigcap_{i \in S} A_i\right) + \mathbb{P}\left(\bigcap_{j \notin S} B_j\right) \leq \mathbb{P}(A_i) + \mathbb{P}(B_j) \\
&\leq (M - |S|) \max_{i \in S, j \notin S} \mathbb{P}(D(\gamma_j \|\hat{\pi}) > D(\gamma_i \|\hat{\pi})) + |S| \max_{i \in S, j \notin S} \mathbb{P}(D(\gamma_i \|\hat{\mu}) > D(\gamma_j \|\hat{\mu})) \\
&\leq (M - |S|)^2 \max_{i \in S, j_1, j_2 \notin S} \mathbb{P}(D(\gamma_{j_1} \|\gamma_{j_2}) > D(\gamma_i \|\gamma_{j_2})) + |S|^2 \max_{i_1, i_2 \in S, j \notin S} \mathbb{P}(D(\gamma_{i_1} \|\gamma_{i_2}) > D(\gamma_j \|\gamma_{i_2})) \quad (50)
\end{aligned}$$

From lemma 1, we know the exponent can be computed as

$$\begin{aligned}
&\lim_{n \rightarrow \infty} -\frac{1}{n} \log \max_{i \in S, j_1, j_2 \notin S} \mathbb{P}(D(\gamma_{j_1} \|\gamma_{j_2}) > D(\gamma_i \|\gamma_{j_2})) \\
&= \min_{i \in S, j_1, j_2 \notin S} \min_{q_1, q_2, q_3 \in C_7} D(q_1 \|\pi_{j_1}) + D(q_2 \|\pi_{j_2}) + D(q_3 \|\mu_i) \triangleq \alpha_7, \\
&C_7 = \{(q_1, q_2, q_3) : D(q_1 \| q_2) > D(q_3 \| q_2)\}, \quad (51)
\end{aligned}$$

$$\begin{aligned}
&\lim_{n \rightarrow \infty} -\frac{1}{n} \log \max_{i_1, i_2 \in S, j \notin S} \mathbb{P}(D(\gamma_{i_1} \|\gamma_{i_2}) > D(\gamma_j \|\gamma_{i_2})) \\
&= \min_{i_1, i_2 \in S, j \notin S} \min_{q_1, q_2, q_3 \in C_8} D(q_1 \|\mu_{i_1}) + D(q_2 \|\mu_{i_2}) + D(q_3 \|\pi_j) \triangleq \alpha_8 \\
&C_8 = \{(q_1, q_2, q_3) : D(q_1 \| q_2) > D(q_3 \| q_2)\}. \quad (52)
\end{aligned}$$

We can upper bound  $\mathbb{P}(F)$  by union bounds,

$$\begin{aligned}
\mathbb{P}(F) &\leq \mathbb{P}(F_1) + \mathbb{P}(F_2) \\
&\leq \mathbb{P}\left(\bigcup_{j \notin S} \{D(\gamma_j \|\hat{\pi}) - D(\gamma_j \|\hat{\mu}) > 0\}\right) + \mathbb{P}\left(\bigcup_{i \in S} \{D(\gamma_i \|\hat{\mu}) - D(\gamma_i \|\hat{\pi}) > 0\}\right) \\
&\leq |S|(M - |S|)^2 \max_{i \in S, j_1, j_2 \notin S} \mathbb{P}(D(\gamma_{j_1} \|\gamma_{j_2}) - D(\gamma_{j_1} \|\gamma_i) > 0) \\
&\quad + |S|^2(M - |S|) \max_{i_1, i_2 \in S, j \notin S} \mathbb{P}(D(\gamma_{i_1} \|\gamma_{i_2}) - D(\gamma_{i_1} \|\gamma_j) > 0). \quad (53)
\end{aligned}$$

From lemma 1, the exponent can be computed as

$$\begin{aligned}
&\lim_{n \rightarrow \infty} -\frac{1}{n} \log \max_{i \in S, j_1, j_2 \notin S} \mathbb{P}(D(\gamma_{j_1} \|\gamma_{j_2}) - D(\gamma_{j_1} \|\gamma_i) > 0) \\
&= \min_{i \in S, j_1, j_2 \notin S} \min_{q_1, q_2, q_3 \in C_9} D(q_1 \|\pi_{j_1}) + D(q_2 \|\pi_{j_2}) + D(q_3 \|\mu_i) \triangleq \alpha_9 \\
&C_9 = \{(q_1, q_2, q_3) : D(q_1 \| q_2) > D(q_1 \| q_3)\}, \quad (54)
\end{aligned}$$

$$\begin{aligned}
&\lim_{n \rightarrow \infty} -\frac{1}{n} \log \max_{i_1, i_2 \in S, j \notin S} \mathbb{P}(D(\gamma_{i_1} \|\gamma_{i_2}) - D(\gamma_{i_1} \|\gamma_j) > 0) \\
&= \min_{i_1, i_2 \in S, j \notin S} \min_{q_1, q_2, q_3 \in C_{10}} D(q_1 \|\mu_{i_1}) + D(q_2 \|\mu_{i_2}) + D(q_3 \|\pi_j) \triangleq \alpha_{10} \\
&C_{10} = \{(q_1, q_2, q_3) : D(q_1 \| q_2) > D(q_1 \| q_3)\}. \quad (55)
\end{aligned}$$

Due to the fact that the objective function of (51), (54) can only be zero at a collection  $q_1 = \pi_{j_1}, q_2 = \pi_{j_2}, q_3 = \mu_i$ , which are not in the constraint sets due to our clustering assumption (1). And the objective function of (52), (55) can only achieve zero when  $q_1 = \mu_{i_1}, q_2 = \mu_{i_2}, q_3 = \pi_j$ , which are not in the constraint sets, either. Thus, we can conclude that  $\alpha_7, \alpha_8, \alpha_9, \alpha_{10} > 0$ . From that  $\lim_{n \rightarrow \infty} \frac{\log M(M-|S|)}{n} = 0$ , we get that

$$\alpha\left(\delta_{c3}^{(1)}, \{\mu_i\}_{i=1}^M, \{\pi_i\}_{i=1}^M\right) = \lim_{n \rightarrow \infty} -\frac{1}{n} \log e\left(\delta_{c3}^{(1)}, \{\mu_i\}_{i=1}^M, \{\pi_i\}_{i=1}^M\right) \geq \min\{\alpha_7, \alpha_8, \alpha_9, \alpha_{10}\}. \quad (56)$$

## F GL test is not exponentially consistent when both typical and outlier distributions forming clusters

Since all the typical and outlier distributions can be distinct, the GL test of replacing the true distribution in (2) by their MLEs leads to identical likelihood estimates for each hypothesis. Thus, the GL approach is not applicable here. One could apply the test in (6) to this problem, but the following example shows that the test in (6) is not universally exponentially consistent, even if condition (1) holds.

As shown in [4], the error exponent of the GL test in (6) is established by showing the following optimization problem has a positive minimum value

$$\begin{aligned} & \min_{q_1, q_2, \dots, q_M \in C_{(S, S')}} \sum_{i \in S} D(q_i \| \mu_i) + \sum_{j \notin S} D(q_j \| \pi_j) \\ C_{(S, S')} = & \left\{ (q_1, q_2, \dots, q_M) : \sum_{i \in S} D\left(q_i \| \frac{\sum_{k \in S} q_k}{|S|}\right) + \sum_{j \notin S} D\left(q_j \| \frac{\sum_{k \notin S} q_k}{M - |S|}\right) \right. \\ & \left. \geq \sum_{i \in S'} D\left(q_i \| \frac{\sum_{k \in S'} q_k}{|S'|}\right) + \sum_{j \notin S'} D\left(q_j \| \frac{\sum_{k \notin S'} q_k}{M - |S'|}\right) \right\}. \end{aligned} \quad (57)$$

We consider the case where  $M = 1000$ ,  $S = \{1, 2\}$ , the typical and outlier distributions are specified by the following

$$\begin{aligned} \mu_1 &= \left(\frac{1}{4}, \frac{1}{2}, \frac{1}{4}\right), & \mu_2 &= \left(\frac{1}{5}, \frac{7}{15}, \frac{1}{3}\right) \\ \pi_3 &= \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right), & \pi_4 = \pi_5 = \dots = \pi_{1000} &= \left(\frac{247}{500}, \frac{32}{125}, \frac{1}{4}\right) \end{aligned} \quad (58)$$

It can be verified the clustering condition (1) holds for this example. However, if we let  $q_1 = \mu_1$ ,  $q_2 = \mu_2$ ,  $q_3 = \pi_3$ ,  $q_4 = \dots = q_{1000} = \pi_4$ ,  $S = \{1, 2\}$  and  $S' = \{1, 2, 3\}$ , then

$$\sum_{i \in S} D\left(q_i \| \frac{\sum_{k \in S} q_k}{|S|}\right) + \sum_{j \notin S} D\left(q_j \| \frac{\sum_{k \notin S} q_k}{M - |S|}\right) \geq \sum_{i \in S'} D\left(q_i \| \frac{\sum_{k \in S'} q_k}{|S'|}\right) + \sum_{j \notin S'} D\left(q_j \| \frac{\sum_{k \notin S'} q_k}{M - |S'|}\right) \quad (59)$$

also holds, i.e.,  $(q_1, q_2, \dots, q_M) \in C_{S, S'}$ , which means the error exponent in (57) is equal to zero. Thus, the test in (6) is not universally exponentially consistent for the case where both typical and outlier distributions form clusters.



## References

- [1] A. Tajer, V.V. Veeravalli, and H.V. Poor, “Outlying sequence detection in large data sets: A data-driven approach,” *IEEE Signal Processing Magazine*, vol. 31, no. 5, pp. 44–56, Sept 2014.
- [2] R. J. Bolton and D. J. Hand, “Statistical fraud detection: A review,” *Statistical science*, pp. 235–249, 2002.
- [3] J. Chamberland and V. V. Veeravalli, “Wireless sensors in distributed detection applications,” *IEEE Signal Processing Magazine*, vol. 24, no. 3, pp. 16–25, 2007.
- [4] Y. Li, S. Nitinawarat, and V. V Veeravalli, “Universal outlier hypothesis testing,” *IEEE Trans. Inform. Theory*, vol. 60, no. 7, pp. 4066–4082, 2014.
- [5] Y. Bu, S. Zou, Y. Liang, and V. V Veeravalli, “Universal outlying sequence detection for continuous observations,” in *Proc. IEEE In. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 4254–4258.
- [6] A. Banerjee, S. Merugu, I. S Dhillon, and J. Ghosh, “Clustering with bregman divergences,” *Journal of machine learning research*, vol. 6, no. Oct, pp. 1705–1749, 2005.
- [7] Y. Li and V. V Veeravalli, “Outlying sequence detection in large datasets: Comparison of universal hypothesis testing and clustering,” in *Proc. IEEE In. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 6180–6184.
- [8] K. Chaudhuri and A. McGregor, “Finding metric structure in information theoretic clustering,” in *COLT*. Citeseer, 2008, vol. 8, p. 10.
- [9] M. R Ackermann, J. Blömer, and C. Sohler, “Clustering for metric and nonmetric distance measures,” *ACM Transactions on Algorithms (TALG)*, vol. 6, no. 4, pp. 59, 2010.
- [10] R. Nock, P. Luosto, and J. Kivinen, “Mixed bregman clustering with approximation guarantees,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2008, pp. 154–169.
- [11] S. Lloyd, “Least squares quantization in pcm,” *IEEE Trans. Inform. Theory*, vol. 28, no. 2, pp. 129–137, 1982.
- [12] M. Blum, R. W Floyd, V. Pratt, R. L Rivest, and R. E Tarjan, “Time bounds for selection,” *Journal of computer and system sciences*, vol. 7, no. 4, pp. 448–461, 1973.
- [13] Y. Li, S. Nitinawarat, and V. V Veeravalli, “Universal sequential outlier hypothesis testing,” in *Proc. IEEE Int. Symp. Information Theory (ISIT)*. IEEE, 2014, pp. 3205–3209.